

【講演1】

Michal Kosinski (スタンフォード大学 リサーチフェロー)

“Prediction of personal attributes from digital records”

Thank you. *Konnichiwa*, good afternoon. It's a great pleasure to be here. I actually arrived a few days earlier so I had a chance to enjoy my time in Tokyo so far. I have never been before, great people, great food, great outdoors. I heard you have a winter as well, but you can't say from now.

Thanks for a very kind introduction. I must admit it's really easy to write popular papers when you write about Facebook and big data and subjects of this kind. Journalists very much like it, so I guess it's much easier to attract public attention when you do not use any equations actually, I just have some screenshots from Facebook.

But actually I do use some equations. Though Ken asked me to promise, I had to promise before I came here that I will not show you any equations so given his presentation, now I feel that it's a bit of a cheating because he gets all the smart equations and I don't have any, so I'm really sorry for that.

What I do? I do study humans in the digital environment. So, I'm not really sure what I should be calling myself. I'm a psychologist by education, but I do actually work at the computer science department, so I guess, I'm mining data, I'm a big social data scientist or a computational social scientist and I try, I focus on understanding humans in the new digital environment. And before I would start my presentation, I will actually ask you for a few things.

First of all if I speak too fast, I very often get excited about the story that I have here, please slow me down. Can we have someone? Can we have a police here, slowing me down police? Tomoya, will you ask me to slow down? Thank you. So, that's the first thing.

Second thing is that there is not so many of us here today, so if you guys have any questions, just shoot immediately, right, like just raise your hand and I'll try to kind of give you voice, let you ask the question. I think it will be much more productive instead of waiting for the end of the presentation. Before I start, I'll actually show you a short 2-minute video here that is a good introduction to the subject. It doesn't work though.

[Video]

So, I'm not sure if you agree with me but this is increasingly the reality we are living in, right? And it's not only the cell phones, it's also computers, all the digital devices that surround us, now the new Samsung cell phone is kind of connected with your watch, so kind of everything you do is now mediated or accommodated by digital devices. And it

kind of brings two issues to my attention. First one is that it gives us a lot of data to study human beings. Everything we do now is being recorded.

You send emails, so kind of do communications while it is recorded, or you kind of talk to them on a cell phone which is also being recorded on our device. I'm not sure if you know guys that your voice conversations, even if they're not recorded, the kind of a voice amplitude would be recorded. So, later GSM network or the scientists can study how your voice changed during conversation and so on. Pictures, GPS location, everything we do is basically recorded which brings a lot of new data to psychologists.

But there is a second issue, that this also changes our environment. Like this girl here, she didn't really look so happy to be surrounded by people who are more of kind of a cyborg, like a combination of a human being and a smart phone, rather than a human being, so it also changes kind of the environment in which we live in and which is also very interesting from a psychological point of view. Is it good for us, is it bad?

I have kind of always tried to be rather neutral so I believe it's just the environment we live in, it's neither good or bad, and if you like it, it's perhaps good, but definitely affects what's going on, affects our kind of psychological relationships with other people and so on, so there are kind of two big issues for psychologists here.

And you can call it big data. I like to call it digital footprints. Basically, those are records of the data collected by digital devices. And I have a little slide here showing a few kinds of data you can kind of get access to here. You have smart phones, other connected devices, kind of traditional laptops and tablets, and you can see how many of the areas of the life, how many life areas are now accommodated by online services and that in effect recorded rightly. The computers know where you are, your cell phone tracks you all the time you walk around.

Actually, your cell phone tracks much more than your physical location. There are also kind of very sensitive sensors inside, for instance accelerometer, so even we can test how quickly, kind of, how jumpy you are, whether you kind of start to walk very quickly or just very calm, again, which could be a proxy for your mood but I'm not sure if you know that if you keep your cell phone on your bed or in your pocket when you sit, the sensor is activated enough to figure out your heartbeat. So it can also record or track your heartbeat, your temperature and many other biological indicators throughout the day.

But there's much more of course. You take pictures of yourself, so it's not only information about you, but it's also information about everyone around you. So, even if you say, oh I do not use any digital devices myself, I'm out of this game, this is not true, because you have friends who take pictures of you and tag you on it and put it on Facebook.

Moreover you quite like to use a credit card which is also recording what you buy, when you buy, why you buy and so on. You have a bus ticket, you come into the subway, I hope you realize there are basically Master card and/or PASMO card here, the data coming is heavily analyzed. You can analyze the cards or behavior, how early you go and start your day, what time do you come back from work, what are the stations you go to and so on.

My data, the data I work with usually comes from social networks, and it's not because I kind of particularly like social networks, though actually I think it's a great device, it's a great thing, but it's just simply that I got access to a lot of data from this source. And actually, I perhaps should mention it before I go any further, the data, the results you can see here, those are based on the data that we share, my research group shares it with other interested researchers. So, if you guys are interested in studying the kinds of things we are studying, you are more than welcome to reach out, send us email.

We have a website called 'mypersonality.org', where you can actually see 'mypersonality.org', where you can actually see what's available. We have data of more than 6 million people. We have their personality scores, IQ scores, happiness scores, morality scores. We have a bunch of psychometric measures applied to them and we also have the records of their Facebook profiles. So they gave us access to the Facebook profile so that we have matched this data in an interesting way, let's say geographical location, aggregate personality. And so there is a whole bunch of things that could be studied and my research team is not large enough to cover everything.

Okay so we have this data, and now the big question is so what? Like you have this data, what can you do with it? Actually one of the questions that I started with is trying to see whether there are any individual differences in digital footprint that make sense? So in psychology—how many psychologists do we have in the room? Can the psychologists raise their hands? Psychologists raise your hands. One, two, three, four, five, six, seven, eight, nine, ten, eleven, okay, so still minority. I was expecting more psychologists. That's okay.

For instance in psychology we distinguish between males and females, between genders, of course we look at age, we look at IQ but we also look at personality and most popular personality model at the moment is so called Big Five, five factors of personality model which distinguishes between two, introverted and extroverted, well organized and spontaneous, emotional and easily stabilized, curious and very stable and what else, agreeable and disagreeable, and I guess that's it, that was five, right?

My first research question, when I started researching was checking whether those psychological dimensions that we kind of believe guide your behavior in the offline world,

whether the same dimensions will guide your behavior in the online world, whether your digital footprint will be related to your psychological traits. And in fact, it turns out that it is quite related. I will show you a few examples, and there actually will be a bit of a quiz later on. So, one thing we did, we looked at language used on social media. So, Facebook status updates and we analyzed this language in the context of psychological traits.

In this case, we looked at happiness, and you can see that happy people, they use words, abeyance, summer, practice, but also, I have no idea why – actually, this is another interesting question I hope maybe, one of you guys will be able to address, happy people, they speak about swimming a lot. There could be some psychological theories to build up to explain that, but this is kind of an example of this becoming a very interesting way of studying people, and answering all the questions, but can also bring you new knowledge, can actually show you some new phenomenon that perhaps we didn't notice before.

Okay now, there is a quiz for you. It is not yet quiz because I kind of have solved it. So those are words that are used by unhappy people, and if you see people using those words, and they used a lot of them, and actually this psychologically makes a lot of sense, like sad people, they don't look very happy. Okay, this is personality trait, one of the big five traits, can you tell me what kind of people would use those words?

Male Speaker

Japanese.

Michal Kosinski

Japanese, yeah, but there is also a psychological trait that is kind of – please note that this is based on American sample. The same words can mean different things in different cultures, so what is a very specific niche thing in America might be a mainstream kind of cultural thing in Japan, for instance... Any guesses, what kind of personality trait it might be? Come on, psychologists, introverted. Internet, computer, sitting at home alone, not really interacting so much with other people, kind of, introverted love internet.

Actually there is an interesting study, I can't remember the reference now, but when the text messaging first showed up, psychologists were like, oh, my god, this will completely destroy the social interactions between people. Now they will just message to each other, they will lose social skills and so on, and that was a basically big issue then. But what actually turned out to be the effect of text messaging is that people who are shy and introverted, those who traditionally had big trouble to start psychological meaningful relationships, because they are basically too shy to approach the others, text messaging gave them this easy, low hanging step towards starting social interactions.

So, I'm too shy, I like this girl a lot but I'm kind of too shy to approach, I'm an introvert, so I kind of send her a message and she sends me a message back and I can actually start conversation, a relationship, something that was not available to the kids in the past, not so easily. So, actually, this is an example of how internet, internet communication can actually help to bridge relationships rather than split them up.

Psychologists, what kind of personality trait that would be? Love, party, excited, weekend, baby, girls, amazing, chilling. Extroverted, of course. Wonderful, family, excited, beautiful, blessed, there must be Jesus somewhere there, church, Thanksgiving, very kind of positive, social kind of oriented agreeable person, high on agreeableness. Next one, hell, darn, not sure if I can say this word aloud, this is from a published paper, so I hope you don't mind, and those are actual words used by people on their status updates, disagreeable people, competitive, working alone, disagreeable, very ambitious and so on.

Shopping, husband, first day, weekend, football, any ideas? Conservative, people high on conservatism. Okay, but it's not only psychological traits that you can look at, so this is actually a demographic trait. What you see here, and what I won't go into it, and those are not just simple word count, so it's not that we simply counted words, we actually used quite sophisticated analysis, latent Dirichlet allocation to look at topics that people talk about, because the same word can mean two different things in two different contexts.

Like jaguar can mean a car but it can also mean an animal. So, this method actually allows you to have the same word in two different contexts kind of together. So those are topics, those are groups of words used very often by a very specific demographic group.

Logic, opinions, opinion more on political society, some swear words here, games, Xbox, playing online, country, riots, freedom, government, democracy, liberty, political country. What kind of demographic group can you see here? Come on guys, this is easy, you must have some guesses. I need a guess from you. What demographic group? Who do you see here? Imagine a person; try to because obviously you all have also psychological intuition in you. What kind of person, if you just see people or a person talking about but you don't know who they are, who do you think would be the person that you see.

Male Speaker

Students.

Michal Kosinski

Students, okay. Male student or a female student? Men or women?

Male Speaker

Men.

Michal Kosinski

Men, exactly, and then you have males here, and now obviously you have a very big difference in language used by females. Puppy, sweetheart, happy, so, so loving, so friend, no politics, no swearing, no shit, no logic, completely different, right? And obviously on average most of the language we use is not different between genders. We use 90 – well I forgot the actual number but a great majority of the language is the same between genders. But have you seen a guy saying, “Oh, I’m so excited”.

It’s not what guys do. And also girls will not talk so much about government, freedom and so on. So it’s kind of easy to stereotype as well, right. You think about stereotypical husband and wife, husband will be interested in politics, women think politics is not really so exciting, so again, I’m not trying to build stereotypes here, but I’m saying that kind of big data, this is based on, I think, it was 300,000 people here, kind of shows that there is actually a very nice signal there.

And of course, you can go much deeper. Again, those are female-dominated topics, and those are male-dominated topics, and you can kind of look at much more in terms of a spectrum of psychological traits. Okay but it’s not only language. The next thing I looked at was interests. People are also interested in different things, and I used, as a proxy for interest, I used Facebook likes. Who here is not using Facebook? Can I see people who do not use Facebook? Come on, there must be someone who does not use Facebook. Not even one person? That’s a big progress actually I’ve never seen an audience where everyone is using Facebook. Congratulations, I think Facebook is a great thing to use.

So, you know what Facebook like is. You actually don’t even have to log in to Facebook, you can just like be on the internet, you visit the website that you like, and you can click ‘Like’ button and this will be saved on your Facebook profile. And now what we can do, we can see, we can look at the differences between different personalities, different psychological profiles, in terms of what they like, what preferences they have. And I will give you an example here.

Two groups, on one psychological trait and on one extreme, you have people who like Mozart, and thunderstorms and science and The Godfather, the movie or the book and they also like Morgan Freeman’s voice. There’s no Morgan Freeman? I also actually liked his voice, so hopefully I’m there. And on the other side you have people who like Harley Davidson and some brands that’s – Tomoya said that he doesn’t recognize any of them, which I think actually speaks very good of you.

Another example, writer, Charles Bukowski, Oscar Wilde, Leonardo Da Vinci, Plato, versus, umm, I don't read. NASCAR, which is basically car racing. This will be open as to experience, open, liberal, progressive, great people here; conservative, traditional people on this side.

Any guesses? This is actually a demographic trait here. Wartime plan versus NOH8 campaign, people, a bit of American culture gets American for example, Ellen DeGeneres is a very famous lesbian American TV presenter, and in fact this is a language, those are preferences differences between straight and gay individuals.

Psychology again, spontaneous, not well organized, not well planning, joining your fans, well, I'm not really very familiar with those objects here but are kind of basically books and well, I mean again versus law officer, national law enforcement, accounting, you have to be highly conscientious to choose a career in accounting. And so on.

I can basically show you quite a few examples, but basically, now the question is, okay so we see those differences, but as actually Ken said, okay, if you have enough data, you will also see some differences but now the question is how meaningful those differences are?

Right so, okay, it kind of depends on signal, how meaningful this signal is, whether kind of it really makes a difference, whether by just looking at the preferences of someone online, or checking the language, can computer automatically actually differentiate between male and female, or a person who is open-minded or conservative? Let me actually again use a video here to cheer you up a bit. And this will be an introduction to the next research that I will present you in a second.

[Video]

Yes, this was basically an introduction to one of the research projects that we have done in the last 2 years. We took 60,000 volunteers from Facebook, we recorded a number of psychological and demographic traits like personality, IQ, but also we checked whether their parents were divorced or not, whether they take drugs, whether they smoke, whether they are gay or straight and we also recorded their digital footprint, and we chose a very simple digital footprint, Facebook likes and our Facebook likes are much less intrusive than your usual digital footprint. Facebook likes is something that is public. You click on like and you know that all of your friends will see, right?

So actually, gay people do not really click on the like I'm gay. Actually a few people do but we have removed them from the sample, as this would be too easy. So, people whose parents were divorced did not click on the like, oh my parents are divorced. I've never

actually seen a thing like that, maybe it exists but I've never kind of encountered it. It's a very public information so people kind of self-censor themselves.

I cannot show you the results but I also looked at the web browsing data, so in Microsoft, Google, Facebook, and other big internet companies, they also collect your web browsing data, they look at the websites you are visiting and here the accuracy is even higher because people, when they browse the internet, they do not feel like someone observes them. Obviously, if you actually, sit down and think about it, obviously your internet provider, your government, Google, Microsoft, they all record this data, like if you're more on internet, it's even more revealing.

But what you will see here is that you can use this very simple kind of signal, Facebook likes and still be very accurate in predicting. And this is one of the advantages, one of kind of the magic things that is being offered by big data. The very simple methods, and very crude, pervasive, generally available piece of data, will give you very accurate and insightful results. So, what we did, I won't hover in for too long, but basically we also used quite simple approach, we didn't use very complicated machinery methods.

What we did is we created a metrics where you have users and the things they liked, and we put 1 in the cells where you had a connection between user and a like, and what we did later, we simply factor-analyzed it, to be honest which is rather singular value decomposition, it was like a PCA kind of a procedure that takes the huge metrics where you have many users and hundreds of thousands of likes, and it kind of flattens it, so we have many uses and only few, in this case 100 components, so it's kind of a simplification of this user metrics. It has many advantages but I'll skip it.

Then what we did, we trained a very simple model. We trained a cross-validated linear regression. It's very difficult to be simpler than that. Linear regression is one of the kind of simplest models you can imagine. So, there cannot be, there was no over-fitting, it was cross-validated, very simple method. So, actually, what I'm saying here is that the results that you'll see in a second are at baseline, but if you use more advanced methods, more intrusive data, perhaps more data, you can only improve on it. So, it's already pretty accurate and I believe it can be significantly improved in the future.

To show you some results, how accurate we were, let's start with some demographic traits. Accuracy here is measured through the area under the curve, which is an equivalent of the probability of being right when you're presented with two people of two different classes. So, let's say we have gender here, this number 0.93 indicates that if you show to a computer, male and female, man and woman, computer in 93% of cases will correctly assign them to a class.

Obviously, I hope that you see that 50 is a baseline, because if I know nothing and you

show me man and woman, randomly I can assign them to classes and I will be correct in 50% of the cases. So, 50% is baseline, 1 is the perfect accuracy. If I always classify them correctly, I'm perfectly accurate, this will be 1.

What you can see here is that you have few things that you can predict very accurately. Gender, 0.93, but also race, 0.95, even higher than gender, which actually from a psychological point of view is pretty interesting because it indicates that the differences in behavior between men and women are smaller than differences in behavior between black American and white American.

Actually from the social point of view, it's kind of a big shot. We have well known gender differences. It seems that races are even more different. You have religion, Christianity versus Islam, Democrat versus Republican, gay versus straight male, 0.88. Relatively low accuracy was achieved, so kind of much closer to the baseline in parents together at 21, so 'yes' parents together, 'no' parents divorced.

But actually, when I entered this variable into analysis, I didn't expect to see any result. I was thinking, come on, how on earth divorce of your parents could change your behavior on Facebook? Why do you like stuff enough to actually allow the computer to distinguish between people whose parents are divorced or not. It turns out that it does.

The kind of divorce of your parents affects you long term, and what you even like on Facebook, music that you listen to, movies that you like, it will affect your preferences. Also, mind you, that this accuracy is also affected by the mistakes in the kind of validation data, in the ground truth. Some parents might have divorced at 22. For us, this person would be in a non-divorce class. Also, some people might have lied about their gender, but still, I would say this introduces additional error, you can see quite a lot of accuracy here.

We also looked at psychological traits and what you see here is Pearson Correlation Coefficient, so correlation between the ground truth, actual personality trait and what was predicted from Facebook likes. And you can see age is pretty easy to predict, you can also predict the properties of the network, so number of friends, and how dense is the network of friendship, intelligence 0.4 prediction, but actually when you control for the fact that the IQ questionnaires are not perfectly accurate. We call it correction for attenuation. So if you dis-attenuate the result, to account for the fact that those personality questionnaires are also not 100% accurate, and you can do it actually using Cronbach's alpha.

From Cronbach's alpha you know that the personality is not accurate so you can control for it, and you can see that the accuracy is actually pretty high. This is comparable with a short personality questioning. So, I can either invite you to my office as a psychologist and ask you to fill this questionnaire or I can just simply say, "Hey, can you

show me your Facebook likes?" I'm kind of achieving very similar prediction accuracy. So I hope that those results show you those differences that I showed you before between men and women, between different personality traits, are actually pretty meaningful. They allow computer to very accurately predict those traits, figure out who you are.

Now we could go even further and say, okay so how this compares with human accuracy? And actually with Youyou here, Youyou is the first author on this paper and we recently ran a study in which we were looking at how accurate are computers at predicting your personality based on likes. We used similar methodology to the one I showed you before.

Actually we used LASSO regression which I strongly recommend to you guys, LASSO regression can take hundreds of thousands of millions of variables very briefly reduced in number, it's a very good prediction mechanism, strongly recommended. And what we found is we were able to compare the accuracy of the computer with the accuracy of the human judge.

We first took personality scores of people, so kind of how they self-reported personalities, what scores they got when they were filling personality questionnaire, and then we asked different other people to judge personality of those individuals. So it's like I'm filling the personality questionnaire and then Tomoya is my friend, and I say, "Hey Tomoya-sensei, can you please fill the personality questionnaire to judge my personality?"

And what you can see here is that the correlation between my own personality score and the personality score judged by my work colleague is only 0.27. So my real personality versus personality assessed by my work colleague is somewhat correlated, 0.3. If you take a friend, or someone who lives with me at one address, this accuracy is even higher, 0.45, 0.44, pretty high. Family member, even higher, family members will know you even better than your friends and quite obviously, spouses, your wife, your husband, your boyfriend, girlfriend, they know you the best so the correlation of their assessment to my assessment is 0.58.

Now what you can do, you can now see how computers kind of automated assessment would perform in this context. So, what we did, we took basically people with different number of likes, because the more likes you have the more I know about you because if you only have two or three likes, I know very little about you. But if you have, like some people have 300, 400, 500 likes, I know more about you so assessment presumably will be more accurate. And in fact, you can see that the computer accuracy rose with number of likes available. So, if you have only like 10 likes or 2 likes, the correlation between computer assessment and my own assessment will be 0.2. If I have 300 likes, it's about 0.6.

Now, let's compare it between – compare kind of computer assessment and friends, and

kind of other people doing assessment, and you can see this is computer expected accuracy. So basically, people on average have around 220 likes so if you give me an average person from Facebook, computer will be accurate 0.56 across different traits. You can see that it's really close to your wife or husband, but really there is not much difference. Statistically, I would say it's actually insignificant there. The thing is that as we collect more data about you, you can see that this line is just growing.

Obviously this is a logistic scale so it's kind of flattening if you have like a non-logistic scale but the more information I have about you the more accurately computer is able to predict your individual traits which actually brings us to a kind of a – for me as a psychologist, I don't want to say kind of life changing, but is basically world-view changing discovery that I was always convinced that humans, and of course especially psychologists, we are so bloody good at judging other people's personalities, we're kind of experts at it, like we study for years, we have tools, we do learn theories and also our brains evolved to study each other's personalities. This was crucial for the survival.

In the past when we lived on the Savannah and you met another human being, it was crucial for you to quickly assess whether they are friendly or not friendly, whether they are extroverted or introverted because this could indicate whether they will be a good teammate or maybe they'll be a danger. So, we've actually evolved to be good at judging other people and there is a lot of evidence showing that, I need only 5 seconds with you or 4 seconds with you to kind of know who you are to a large extent.

I'm sure you've kind of heard of all those papers showing that short impressions are very strong signal for a human being. And now what turns out is that computers based on the very simple kind of a data are better than human beings, so kind of approaching spouses on average and you can see this line just goes forever. But they are actually better than human beings at judging very human trends.

And then again, question of so what? What does it mean for the future of science of humanity and so on? And I think that one of the analogies that I particularly like is between what we now observe happens in social sciences, big data social sciences and what happened to the weather science. I'm sure guys you're familiar with weather forecast. I have it on my cell phone, it's actually increasingly accurate but people usually do not realize that weather science, starting with how clouds move, how cold air moves from one place to the other, this is very young science. This is only 100-150 years old. Before that, people believed that science is local – not science, weather is local.

You are living in Tokyo and you're seeing that the sky is red and you're thinking, oh when the sky is red, it will rain tomorrow, or it will not rain tomorrow, whatnot. And I'm sure you actually know many of those kinds of pieces of knowledge, kind of folk

knowledge are still there. Actually, your grandparents or parents study all these like birds' flight, oh then tomorrow weather will be such and such. But we know today that weather is not local. That what we'll see here, whether it rains here today or not, actually depends on the wind patterns of how the temperatures vary, what are the differences in pressure, what is high pressure, low pressure and whatnot.

But people did not realize it for most of our existence, and please note that weather is so crucial for our survival. Farming, hunting, so many things depend on the weather. But it was impossible for human beings to understand the weather simply because they did not have big data. The only data they had was from here now, today. I today see red sky, tomorrow it will rain or not. And actually, modern weather started with the discovery, with implementation of discovery of a telegraph.

Farmers in United States, when they were building a railway network, they had also the telegraph along the railway network. So if it was raining at one place, one telegraph guy, who was sitting there bored, would be like "Oh, nothing is happening in my city today, but it's raining here today". He would like send it to a person who was 100 kilometers away or 150 kilometers away and they kind of noticed that if it's raining here, and the wind is from there, it will rain there.

This is so obvious for us, but they had no idea 150 years ago, like they basically had no idea and then farmers actually were like it was the most important thing at the time that was done across. Farmers realized that this actually is useful and they started studying it, and moving a new science, new understanding. Something that is so obvious for us today, that weather actually is a pattern, it's a global pattern. And I think that psychology is today at the stage where weather science was 150 years ago.

We know a lot about human beings. We know how we behave, we see people, we a lot of local patterns, also we have a local folk knowledge. We have a folk knowledge that what people do when they come into a new environment, they see other people, they have a kind of a lot of acquired folk knowledge, what people should be doing, what people usually do. And I'm thinking that big data, on a macro level, so big social data but also on the brain level so big neuroscience data would bring completely new understanding to how people operate, what drives our behavior and so on. Basically, it will be a same revolution that happened to the weather 150 years ago.

Which covers first point here, but I also think there are many more opportunities at hand like kind of small winds, low hanging fruits that we can get from the methods that I've presented today. For instance, you can imagine that those methods that are really through the computer to judge personality and other traits can be used to cheaply and superfast and quickly run psychological assessment. I'll give you an example that is very

simple, yet very powerful.

Today when you want to hire people, or you want to hire people for the university, you invite people, “Hey, send me application”, and you invite them to take the test, and then you kind of accept them or reject them. But now imagine if I had Facebook and I had 1.3 million people on Facebook, what I can actually do, I can say, I will now analyze everyone on Facebook, I will select 100 people that I believe actually will match my course, and instead of kind of fishing for people, I could just invite those 100 people and say, “Hey guys, you will be great for this job. Would you like to come here?” Or “You would be great for this university for this course; would you also like to come here?” So, kind of job matching, university course matching, maybe dating matching could be greatly improved by computers being able to automatically run their assessment.

And obviously we get a lot of interest, so a lot of marketing companies contacted us because they believe and actually, I agree with that, that by being able to very quickly and automatically understand human beings, understand their kind of psychological profile, you can also improve your marketing targeting and marketing communication. You can improve your messages; improve the products you’re offering to people.

But obviously, there are also risks associated with this technology, and I’m sure, well, I hope at least that some of the results we showed you here were a bit scary. For instance, for me, one of the foundations of the liberal democratic society is that I can share with you the data on my sexual orientation, I can tell you whether I’m gay or straight, or I can tell you what’s my political views, or I can tell you whether I’m a catholic or an atheist, if I want to but if I don’t want to, I don’t have to tell you. And I think this is very important for our security and for our freedom that we have this right to withdraw the information.

But what happens here is that looking at those results you cannot, any longer, withdraw this information from the others. All of your Facebook friends can use a very simple online software to scan your Facebook likes for instance, or Facebook status updates or your email that you exchange with them, and the software will tell with very high accuracy, Tomoya is a liberal person, or is a conservative person.

So something, that before was your private information and was protected by law, I’m not sure about Japan but let’s say in Europe, information about your sexual orientation, about your IQ, about your palette, about your political views, it’s highly protected by law. You cannot record it without my permission and you have to destroy that in sometime and you have to store it in a secure place and so on.

But now what happens and actually it is the case now, and this is not my opinion, it’s the opinion of European Parliament; the internet is illegal in Europe now. According to the law, you cannot have a system running that reveals your IQ to anyone who’d like to

see it, or your sexual orientation for anyone who'd like to see it, so obviously, Europe is not switching the internet off, but believe me there's a lot of thinking going on in European Commission, European Parliament, what to do with this problem? We believe it's important to protect the information of individuals but now if you disseminate it on the internet, that kind of goes against it.

One more thing is quickly assessment, so I said, it's possible to assess people at large scale but maybe once we know that, we're satisfied, and then if someone, even though if I know, that someone scans my Facebook profile to just check out my IQ for instance, I might feel I'm – so I think there's a lot of space for research to understand what people would like to set foot on, and also for policy making, for good policy making, to protect the people from companies evil, or good companies doing stuff that we don't want to do. Overall, I personally actually believe that this technology is like many other technologies, so big social data is a bit like a knife.

We can use knife to kill people, but we can also use knife to cut meat and prepare food and so on, so it can be applied both for good and for bad. And to give you an example, of course, we could have a separate presentation for this bit, but I'll give you an example. In my opinion people love to be assessed, they love to share the data, as long as they have control over it. So, the sample of 6 million people that I told you about, those were volunteers, they donated data to me and they said, please use it, because I gave them choice.

I said, you can participate in my research, you can get the results but you can also help me by collecting the data, and people were very happy and even when we published those results, I got a lot of emails saying all this had amazing result, and I got no email saying this is really creepy. But actually, when you think about it, this is pretty creepy, the kind of results that I've shown you, but because we gave people choice they didn't mind.

And I believe the same relates to assessment. If, let's say, people get to choose, say, I want to participate in this system, I want to be automatically assessed, people will love it. If you just do it without asking, people will hate it, and both in science, participants, and I agree with you and in the corporate world or the government world, you can actually spook out individuals.

This is what I've prepared for you today and there were no questions during the presentation but I hope you have some questions now.

Questioner

Last time we talked about this, some of your studies included more than 50,000 participants, I don't remember exactly but The Daily Mail news article seems very negative about your research. Would you tell me how do you recruit the volunteers for your study?

Michal Kosinski

Okay, that's a very good question, thank you for that. So, what we did, actually, it was all started by my friend and collaborator Dr. David Stillwell, when he finished his undergraduate course, and was just going to do his Masters, he had holidays for which he didn't plan anything. He is an introvert so he likes sitting at home, so what he did and he just finished his undergraduate degree in psychology, so he took an open source personality test, IPIP 100-item proxy for NEO-PI-R, he took this personality test and he put it online and just the timing was 2006, 2007 – when Facebook opened the so-called application API, so everyone could publish like an application on Facebook.

So, he decided to publish personality questionnaire, and he thought, he imagined that maybe 100 people would take it. Maybe 200 people would take it. He sent it to his friends, he didn't use any marketing, he just put it on Facebook, and it turned out that there were up to 1 million people a month, coming to the application, taking the test. It didn't look great. It's actually now switched off, it didn't look great, it was just psychology student made the application so you can imagine it was like pretty simple.

And then we started adding additional tests, IQ, happiness, we had, I think, in total 40 different questionnaires and people just loved it. So they kind of could come, take the test, they could post the results on their Facebook – newsfeed Facebook wall, so then others could see and say, oh this is interesting, I also want to take the test, and in this way, I guess, in psychology it's called 'snowball approach'. In this way we kind of snowballed, I think a total sample of 8 million people, 8 million people agreed for us to use their data, and actually even more people took the questionnaire.

Interesting thing, I don't think it will be possible anymore to copy exactly this solution simply because there is novelty factor involved. So, the first thing usually takes off, all the copies are kind of not so popular because the people already took the questionnaire, those that they were interested. Moreover, now Facebook is much more commercial. At the beginning, Facebook wanted apps to be there, maybe you remember like Mafia Wars, so you had lots of this kind of spam on your wall, one of the spam was our study. Nowadays, we have to pay a dollar to actually before your study to be shown on people's walls.

Questioner

Changed your strategy recently, first it was kind of accidental but now Facebook is more protected, right?

Michal Kosinski

Yes, we try to find new ways of making it interesting and actually if I have internet here – I don't, but you can check yourself. There is a website called youarewhatyoulike.com. When we presented this study, people were very interested. They were like, oh, I'm wondering what will be my personality based on my likes. So, we created this website where you can come, you can click, I want to assess my personality and the website will look at your likes and give you a personality profile, and again, you can donate these likes to research. We do not really have a personality questionnaire there but we get a lot of profile information. And so you can find new wild ways of collecting the data.

And you'll be surprised because I have a traditional psychological admission, so my first study was based on 20 people that I collected in a month, just spend a lot of time collecting variable data, and I think that to create study like that online study, we have to put a bit more work, we have to kind of design something on a computer, it requires a bit more work, but the return is insane.

Because at the end of the day, to create an online questionnaire versus an offline questionnaire is just maybe two times more expensive but the sample size and the amount of data that you get back is not two times more, it's million times more. So, I think it's really worth it. I would encourage all of you guys to go and do study, at least collect the data online.

And there is one more important thing here, I'm not sure how many of you guys are students, but believe me, majority of psychologists and social scientists, they know nothing about programming, they know nothing about computational methods and so on. So, if you also know nothing about programming and computational methods you have to compete with all of those other people.

But if you actually learn, and it's not so difficult, I'm not a computer scientist, if you actually learn them, again it's not so difficult at the end of the day, you have super competitive advantage. It's much more easy to find jobs, it's much more easy to run interesting studies because you don't have to compete with all those other people doing studies in your area, you can have a new area that you created for yourself.

Kensuke Okada

Any other questions?

Questioner

Thank you for the great presentation and I'm interested in the change of the psychological traits, so I feel your research is selective or definitive of the psychological traits, so make just like time-series analysis, or if you have an idea about that.

Michal Kosinski

Thank you for this question, it's really good. Well, I'm surprised that you came up with this idea so quickly. It took me some time to arrive at a simple and elegant solution, and actually, I'm currently working on basically very similar analysis over time. So, you can see how traits are being expressed in different ways across that. But I think, actually, there is even more important issue that we can try to explore, using big social data. Which namely, I'm not sure, if you guys are from media, how Big Five personality was first developed, so that the underlying idea behind Big Five personality is actually pretty elegant.

And this idea is as following. If there are some important individual differences in people, like the people are different in an important way, psychologists, who have kind of developed this personality system or framework, they said, if there are those differences, they will be expressed in language, because over millennia, if people realized it's important that some people are nice and some people are not nice, that will basically have wars starting.

So, what they did, they did linguistic analysis, so they collected all the adjectives described in the book, and then they ran a number of factor analysis so, let's say, they took a person and said, hey, here you have a list of adjectives, can you describe your friend? And you are describing your friend with a list of those adjectives and then some people were looking and saying, open minded, goes along, creative, and it also goes together with liberal. So, those used to be clubbed together. But then they run factor analysis, and they said, all those adjectives, they actually describe one trait, open-mindedness, why did they call it open-mindedness.

And now the thing is that very similar elegant idea can be applied to big social data. In the past, the important individual differences would be expressed in language. Now, what I'm saying is that important individual differences would be expressed in your behavior that now we can record on the human scale.

What I can do actually – I can't say that I know the details of how I'm going to do it, but this is what I'm working towards, is trying to redevelop personality model, IQ model and other kind of psychological dimensions looking at big data only. So, I say, okay, I've showed already this is related to personality. But now let's just look only at the behavior

and now look at dimensions behind it and say, hey, maybe this will be a better description of who we are and how we differ between each other.

Kensuke Okada

I think you may have a lot of other questions but the time comes. I just want to add one thing. He recommended LASSO regression and I believe that LASSO regression can be most easy understanding in Bayesian sense so I would also recommend using LASSO regression and use the Bayesian interpretation. By the way, let us give a big hand to Michal, thank you very much.

Michal Kosinski

Thank you.